

FLP

Theorem

There is **no** deterministic protocol that solves Consensus in a message-passing asynchronous system in which at most one process may fail by crashing

Around FLP in 80 Slides

How can one get around FLP?

Weaken the problem

👁 Weaken termination

- use randomization to terminate with arbitrarily high probability
- guarantee termination only during periods of synchrony

👁 Weaken agreement

□ ϵ - agreement

- ▶ real-valued inputs and outputs
- ▶ agreement within real-valued small positive tolerance ϵ

□ k-set agreement

- ▶ **Agreement**: In any execution, there is a subset W of the set of input values, $|W| = k$, s.t. all decision values are in W
- ▶ **Validity**: In any execution, any decision value for any process is the input value of some process

How can one get around FLP?

Constrain input values

- Characterize the set of input values for which agreement is possible

Strengthen the system model

- Introduce **failure detectors** to distinguish between crashed processes and very slow processes

Paxos

The Part-Time Parliament

- Parliament determines laws by passing sequence of numbered decrees
- Direct democracy: Citizens/Legislators leave and enter the chamber at arbitrary times
- No centralized records: each legislator carries a ledger recording the approved decrees



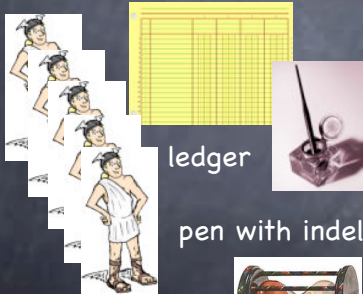
Government 101

- No two ledgers contain contradictory information
- If a majority of legislators are in the Chamber and no one enters or leaves the Chamber for a sufficiently long time, then
 - any decree proposed by a legislator is eventually passed
 - any passed decree appears on the ledger of every legislator

"In a world..."

Political intrigue!

Funky equipment!



ledger

pen with indelible ink

messengers!



hourglass



Λαμπσων
Δίκστρα
Λισκωφ
Mysterious
characters

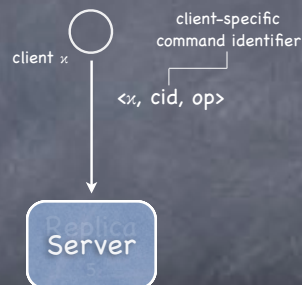
Σθνειδερ Σκεεν
Πννελι Δωλεφ
Δφωρκ



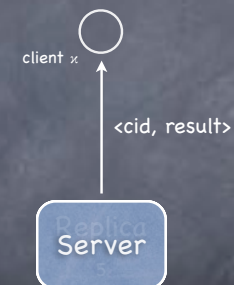
Back to the future

- 👁 A protocol for state machine replication in an asynchronous environment that admits crash failures (key ideas already present in earlier work on Viewstamped Replication by Oki and Liskov)
- 👁 Messages:
 - ❑ between correct endpoints are eventually received
 - ❑ can be lost and duplicated, but not corrupted

The big picture



The big picture



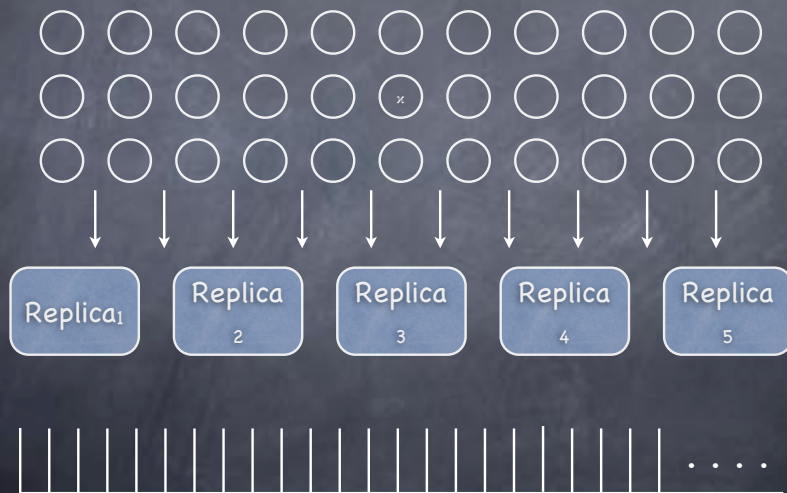
The big picture



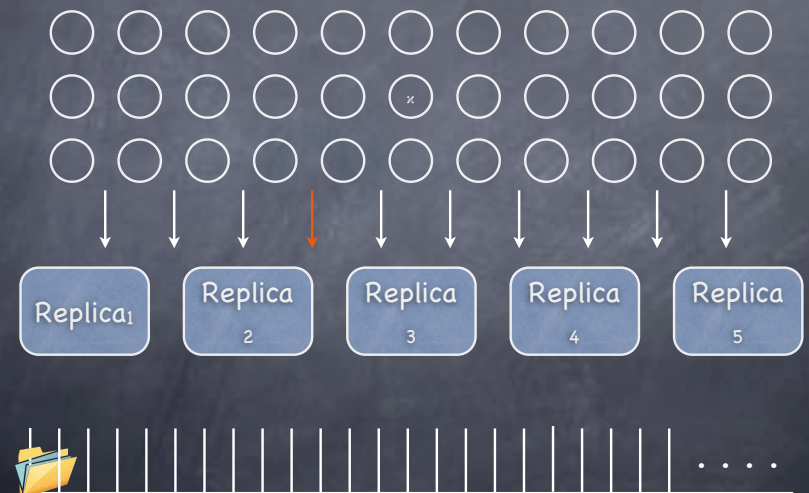
The big picture



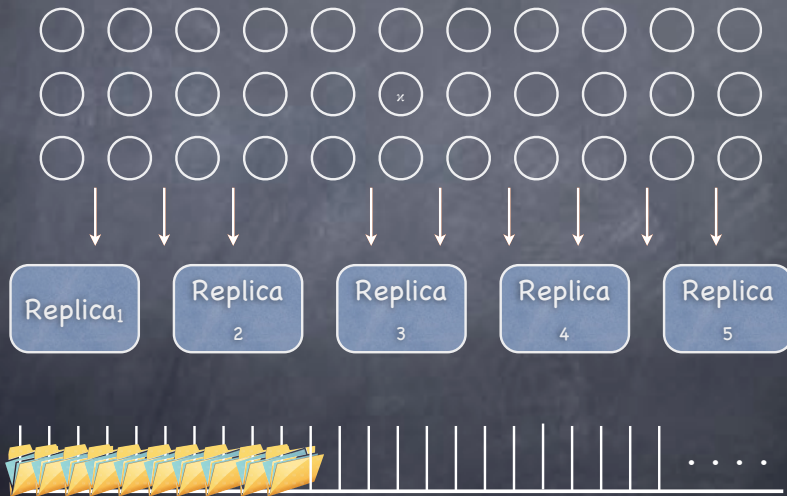
The big picture



The big picture



The big picture



Replicas

- 👁 Receive client requests
 - 👁 Propose command for lowest unused slot to **leaders**
 - 👁 Upon decision, execute commands in slot order
 - 👁 Return result to clients
 - 👁 **Not necessarily identical at any time!**
- 👁 Each replica ρ maintains four variables:
 - ❑ $\rho.state$: the application state
 - ❑ $\rho.slot_num$: next slot for which ρ does not know a decision
 - ❑ $\rho.proposals$: set of $\langle \text{slot number}, \text{command} \rangle$ pairs for past proposals
 - ❑ $\rho.decisions$: set of $\langle \text{slot number}, \text{command} \rangle$ pairs for decided slots

```

process Replica(leaders, initial_state)
var state := initial_state, slot_num := 1, proposals := {}, decisions := {}
for ever
  switch receive()
  case <request, p> :
    propose(p);
  case <decision, s, p> :
    decisions := decisions ∪ {(s, p)}
    while ∃ p' : (slot_num, p') ∈ decisions do
      if ∃ p'' : (slot_num, p'') ∈ proposals ∧ p'' ≠ p' then
        propose(p'');
      end if
      perform(p');
    end while
  end switch
end for
end process
    
```

Replica

```

function perform((κ, cid, op))
  if ∃ s : s < slot_num ∧
    (s, (κ, cid, op)) ∈ decisions then
    slot_num := slot_num + 1
  else
    (next, result) := op(state);
    atomic
      state := next;
      slot_num := slot_num + 1
    end atomic
    send(κ, (response, cid, result));
  end if
end function
    
```

```

function propose(p)
  if ∄ s : (s, p) ∈ decisions then
    s' := min{s | s ∈ ℕ+ ∧ ∄ p' : (s, p') ∈ proposal ∪ decisions};
    proposals := proposals ∪ {(s', p)};
    ∀ λ ∈ leaders : send(λ, (propose, s', p));
  end if
end function
    
```

R4. For each ρ , the variable $\rho.slot_num$ never decreases

```

process Replica(leaders, initial_state)
var state := initial_state, slot_num := 1, proposals := {}, decisions := {}
for ever
  switch receive()
  case <request, p> :
    propose(p);
  case <decision, s, p> :
    decisions := decisions ∪ {(s, p)}
    while ∃ p' : (slot_num, p') ∈ decisions do
      if ∃ p'' : (slot_num, p'') ∈ proposals ∧ p'' ≠ p' then
        propose(p'');
      end if
      perform(p');
    end while
  end switch
end for
end process
    
```

Replica

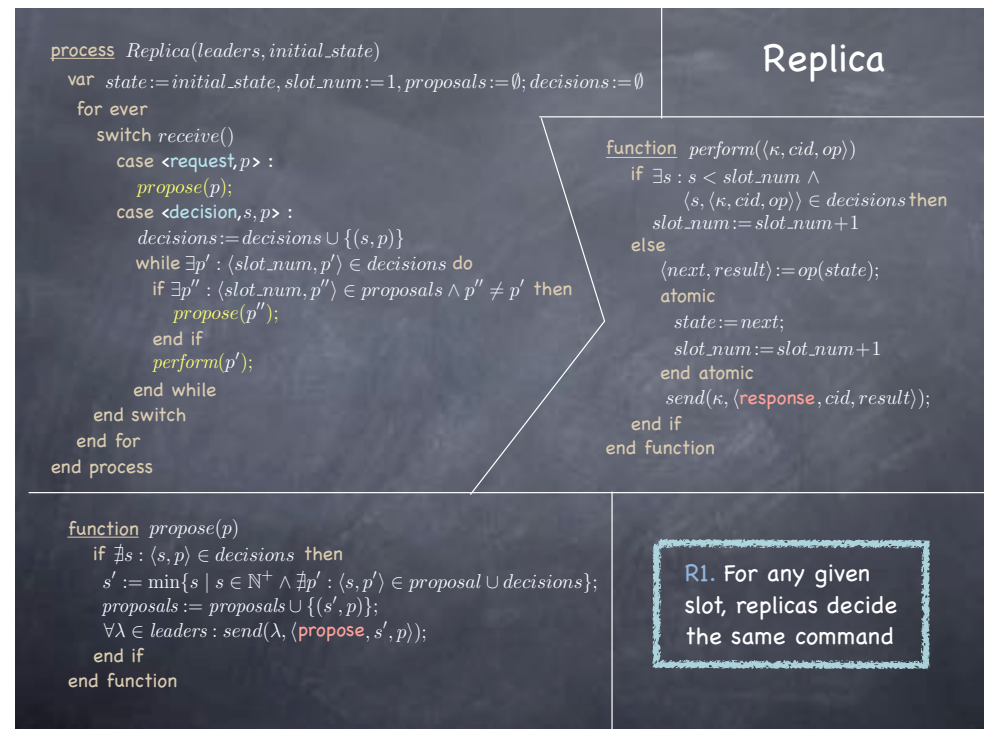
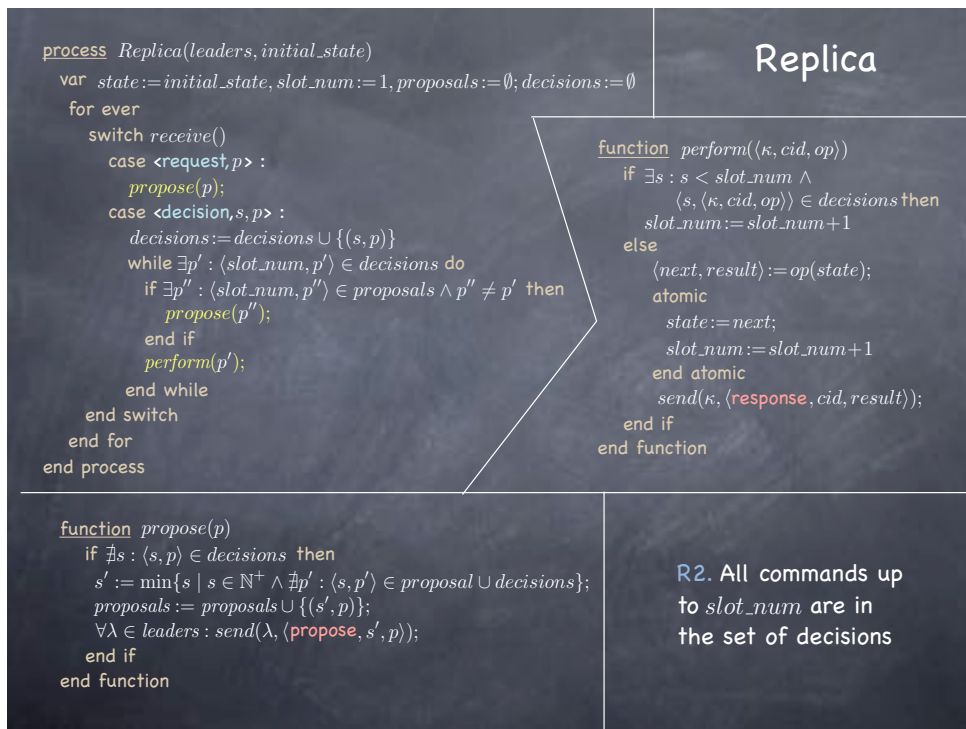
```

function perform((κ, cid, op))
  if ∃ s : s < slot_num ∧
    (s, (κ, cid, op)) ∈ decisions then
    slot_num := slot_num + 1
  else
    (next, result) := op(state);
    atomic
      state := next;
      slot_num := slot_num + 1
    end atomic
    send(κ, (response, cid, result));
  end if
end function
    
```

```

function propose(p)
  if ∄ s : (s, p) ∈ decisions then
    s' := min{s | s ∈ ℕ+ ∧ ∄ p' : (s, p') ∈ proposal ∪ decisions};
    proposals := proposals ∪ {(s', p)};
    ∀ λ ∈ leaders : send(λ, (propose, s', p));
  end if
end function
    
```

R3. For all replicas ρ , $\rho.state$ is the result of applying the operations in $\rho.decisions$ to the initial state in increasing order of slot number s :
 $1 \leq s < slot_num$



Paxos: the basic idea

- Leaders compete to create a permanent mapping between slot numbers and proposals
- The mapping is recorded in "Paxos memory" at a set of state machines called acceptors

Note: Though state machines, acceptors are NOT replicas of each other!
- A leader never proposes a map that may conflict with what is stored in Paxos memory
- A leader, before attempting to create a new map between a slot number for which it knows not a decision and a proposal, "reads" the Paxos memory to check whether such map **may** already exist
- Once a leader learns that a new mapping has become permanent, it informs the replicas

Ballots

- Each leader has an infinite supply of ballots



- The set of ballots of different leaders are disjoint

Ballots

- Each leader has an infinite supply of ballots



- The set of ballots of different leaders are disjoint
- Ballots are lexicographically ordered pairs $\langle seq_no, LId \rangle$

Acceptors

- Send messages only when prompted
- Can crash...
- ...but we assume no more than a minority will
- Need at least $2f+1$ acceptors to tolerate faults
- Each acceptor α maintains two variables:
 - $\alpha.ballot_num$, initially \perp
 - $\alpha.accepted$, a set of pvalues, initially empty
- A pvalue $e = \langle b, s, p \rangle$ is a triple
 - b : ballot number
 - s : slot number
 - p : a proposal
- α **accepts** $e \equiv e \in \alpha.accepted$
- α **adopts** $b \equiv \alpha.ballot_num := b$

A mapping is forever...

- ...once it is accepted by a majority of acceptors – it is then **chosen**
- α accepts a pvalue only if it includes the ballot most recently adopted by α
- To make mapping $\langle s, p \rangle$ permanent, λ needs a majority of acceptors to adopt the ballot of the pvalue that contains $\langle s, p \rangle$

```

process Acceptor()
var ballot_num :=  $\perp$ , accepted :=  $\emptyset$ ;
for ever
switch receive();
case <p1a,  $\lambda$ , b > :
if b > ballot_num then
ballot_num := b;
end if
send( $\lambda$ , <p1b, self(), ballot_num, accepted>);
case <p2a,  $\lambda$ , <b, s, p> > :
if b  $\geq$  ballot_num then
ballot_num := b;
accepted := accepted  $\cup$  {<b, s, p>}
end if
send( $\lambda$ , <p2b, self(), ballot_num>);
end switch
end for
end process
    
```

Acceptor

- On receiving <p1a, λ , b >
 - adopts b iff larger than $ballot_num$
 - returns to λ all accepted pvalues
- On receiving <p2a, λ , <b, s, p>
 - adopts b iff larger than $ballot_num$
 - accepts e if b equal to $ballot_num$
 - returns to λ the current $ballot_num$

Invariants

A1. An acceptor can only adopt strictly increasing ballot numbers

A2. An acceptor can only accept $\langle b, s, p \rangle$ if $b = ballot_num$

A3. An acceptor α can not remove entries from $\alpha.accepted$


```

process Acceptor()
var ballot_num :=  $\perp$ , accepted :=  $\emptyset$ ;
for ever
switch receive();
case <p1a,  $\lambda$ , b > :
if b > ballot_num then
ballot_num := b;
end if
send( $\lambda$ , <p1b, self(), ballot_num, accepted>);
case <p2a,  $\lambda$ , (b, s, p)> :
if b  $\geq$  ballot_num then
ballot_num := b;
accepted := accepted  $\cup$  {(b, s, p)}
end if
send( $\lambda$ , <p2b, self(), ballot_num>);
end switch
end for
end process

```

Acceptor

- 👁 On receiving <p1a, λ , b >
 - ❑ adopts b iff larger than ballot_num
 - ❑ returns to λ all accepted pvalues
- 👁 On receiving <p2a, λ , (b, s, p)>
 - ❑ adopts b iff larger than ballot_num
 - ❑ accepts e if b equal to ballot_num
 - ❑ returns to λ the current ballot_num

Invariants

A4. For a given b and s, at most one proposal can be under consideration by the acceptors: $\langle b, s, p \rangle \in \alpha.\text{accepted} \wedge \langle b, s, p' \rangle \in \alpha'.\text{accepted} \implies p = p'$

A5. Suppose a majority of acceptors has $\langle b, s, p \rangle \in \alpha.\text{accepted}$. If $b' > b$ and $\langle b', s, p' \rangle \in \alpha'.\text{accepted}$, then $p = p'$



Commander

- 👁 A leader holding $\text{ballot_num} = b$ and trying to map slot s to proposal p spawns a new commander thread for $\langle b, s, p \rangle$
- 👁 A commander's mission has two possible outcomes:
 - ❑ **success:** replicas learn that the proposed mapping has been permanently established
 - ❑ **failure:** the leader learns that b may no longer be acceptable to a majority of acceptors



Commander invariants

C1. For any b and s, at most one commander is spawned



A4. For a given b and s, at most one proposal can be under consideration by the acceptors



Commander invariants

C2. Suppose a majority of acceptors has $\langle b, s, p \rangle \in \alpha.\text{accepted}$. If a commander is spawned for $\langle b', s, p' \rangle$: $b' > b$, then $p = p'$



A5. Suppose a majority of acceptors has $\langle b, s, p \rangle \in \alpha.\text{accepted}$. If $b' > b$ and $\langle b', s, p' \rangle \in \alpha'.\text{accepted}$, then $p = p'$


```

process Commander( $\lambda$ , acceptors, replicas,  $\langle b, s, p \rangle$ )
var waitfor := acceptors, pvalues :=  $\emptyset$ 
 $\forall \alpha \in \text{acceptors} : \text{send}(\alpha, \langle \text{p2a}, \text{self}(), \langle b, s, p \rangle \rangle);$ 
for ever
  switch receive();
  case  $\langle \text{p2b}, \alpha, b' \rangle$  :
    if  $b' = b$  then
      waitfor := waitfor -  $\{\alpha\}$ ;
      if  $|\text{waitfor}| < |\text{acceptors}|/2$  then
         $\forall \rho \in \text{replicas} :$ 
          send( $\rho$ ,  $\langle \text{decision}, s, p \rangle$ );
          exit();
        end if;
      else
        send( $\lambda$ ,  $\langle \text{preempted}, b' \rangle$ );
        exit();
      end if;
    end switch
  end for
end process

```



Commander

Must enforce

R1. For any given slot, replicas decide the same command



A5. Suppose a majority of acceptors has $\langle b, s, p \rangle \in \alpha.\text{accepted}$. If $b' > b$ and $\langle b', s, p' \rangle \in \alpha'.\text{accepted}$, then $p = p'$



C2. Suppose a majority of acceptors has $\langle b, s, p \rangle \in \alpha.\text{accepted}$. If a commander is spawned for $\langle b', s, p' \rangle : b' > b$, then $p = p'$

```

process Commander( $\lambda$ , acceptors, replicas,  $\langle b, s, p \rangle$ )
var waitfor := acceptors, pvalues :=  $\emptyset$ 
 $\forall \alpha \in \text{acceptors} : \text{send}(\alpha, \langle \text{p2a}, \text{self}(), \langle b, s, p \rangle \rangle);$ 
for ever
  switch receive();
  case  $\langle \text{p2b}, \alpha, b' \rangle$  :
    if  $b' = b$  then
      waitfor := waitfor -  $\{\alpha\}$ ;
      if  $|\text{waitfor}| < |\text{acceptors}|/2$  then
         $\forall \rho \in \text{replicas} :$ 
          send( $\rho$ ,  $\langle \text{decision}, s, p \rangle$ );
          exit();
        end if;
      else
        send( $\lambda$ ,  $\langle \text{preempted}, b' \rangle$ );
        exit();
      end if;
    end switch
  end for
end process

```



Commander

A higher ballot b' is active: a majority of acceptors may no longer be willing to accept b

```

process Commander( $\lambda$ , acceptors, replicas,  $\langle b, s, p \rangle$ )
var waitfor := acceptors, pvalues :=  $\emptyset$ 
 $\forall \alpha \in \text{acceptors} : \text{send}(\alpha, \langle \text{p2a}, \text{self}(), \langle b, s, p \rangle \rangle);$ 
for ever
  switch receive();
  case  $\langle \text{p2b}, \alpha, b' \rangle$  :
    if  $b' = b$  then
      waitfor := waitfor -  $\{\alpha\}$ ;
      if  $|\text{waitfor}| < |\text{acceptors}|/2$  then
         $\forall \rho \in \text{replicas} :$ 
          send( $\rho$ ,  $\langle \text{decision}, s, p \rangle$ );
          exit();
        end if;
      else
        send( $\lambda$ ,  $\langle \text{preempted}, b' \rangle$ );
        exit();
      end if;
    end switch
  end for
end process

```



Commander

Notify the leader and exit



scout

- Before spawning commanders for ballot b , leader invokes a scout
- Scouts read the Paxos memory to help leaders propose mappings that satisfy **C2**.
- A scout's mission has two possible outcomes:
 - success**: the leader learns that the proposed ballot has been adopted by a majority of acceptors and receives all pvalues accepted by that majority
 - failure**: the leader learns that b may no longer be acceptable to a majority of acceptors

```

process Scout( $\lambda$ , acceptors,  $b$ )
var waitfor := acceptors, pvalues :=  $\emptyset$ 
 $\forall \alpha \in \text{acceptors} : \text{send}(\alpha, \langle \text{p1a}, \text{self}(), b \rangle);$ 
for ever
  switch receive();
  case  $\langle \text{p1b}, \alpha, b', r \rangle :$ 
    if  $b' = b$  then
       $pvalues := pvalues \cup r;$ 
       $waitfor := waitfor - \{\alpha\};$ 
      if  $|waitfor| < |\text{acceptors}|/2$  then
         $\text{send}(\lambda, \langle \text{adopted}, b, pvalues \rangle);$ 
         $\text{exit}();$ 
      end if;
    else
       $\text{send}(\lambda, \langle \text{preempted}, b' \rangle);$ 
       $\text{exit}();$ 
    end if
  end switch
end for
end process

```

Scout

Scout

- gets a majority of acceptors to adopt b
- collects all pvalues that acceptors have accepted while adopting ballots no larger than b

```

process Scout( $\lambda$ , acceptors,  $b$ )
var waitfor := acceptors, pvalues :=  $\emptyset$ 
 $\forall \alpha \in \text{acceptors} : \text{send}(\alpha, \langle \text{p1a}, \text{self}(), b \rangle);$ 
for ever
  switch receive();
  case  $\langle \text{p1b}, \alpha, b', r \rangle :$ 
    if  $b' = b$  then
       $pvalues := pvalues \cup r;$ 
       $waitfor := waitfor - \{\alpha\};$ 
      if  $|waitfor| < |\text{acceptors}|/2$  then
         $\text{send}(\lambda, \langle \text{adopted}, b, pvalues \rangle);$ 
         $\text{exit}();$ 
      end if;
    else
       $\text{send}(\lambda, \langle \text{preempted}, b' \rangle);$ 
       $\text{exit}();$ 
    end if
  end switch
end for
end process

```

Scout

A higher ballot b' is active: a majority of acceptors may no longer be willing to accept b

```

process Scout( $\lambda$ , acceptors,  $b$ )
var waitfor := acceptors, pvalues :=  $\emptyset$ 
 $\forall \alpha \in \text{acceptors} : \text{send}(\alpha, \langle \text{p1a}, \text{self}(), b \rangle);$ 
for ever
  switch receive();
  case  $\langle \text{p1b}, \alpha, b', r \rangle :$ 
    if  $b' = b$  then
       $pvalues := pvalues \cup r;$ 
       $waitfor := waitfor - \{\alpha\};$ 
      if  $|waitfor| < |\text{acceptors}|/2$  then
         $\text{send}(\lambda, \langle \text{adopted}, b, pvalues \rangle);$ 
         $\text{exit}();$ 
      end if;
    else
       $\text{send}(\lambda, \langle \text{preempted}, b' \rangle);$ 
       $\text{exit}();$ 
    end if
  end switch
end for
end process

```

Scout

Notify the leader and exit



Leader

- 👁 Spawns a scout for initial ballot number
- Enters a loop waiting for one of three messages:
 - $\langle \text{propose}, s, p \rangle$ from a replica
 - $\langle \text{adopted}, \text{ballot_num}, pvals \rangle$ from a scout
 - $\langle \text{preempted}, \langle r', \lambda' \rangle \rangle$ from a commander or a scout
- 👁 Each leader λ maintains three variables:
 - $\lambda.\text{ballot_num}$, initially 0
 - $\lambda.\text{active}$, boolean, initially false
 - $\lambda.\text{proposals}$, an initially empty map $\langle \text{slot_number}, \text{proposal} \rangle$
- 👁 Leader moves between **active** and **passive** mode
 - in passive mode is waiting for $\langle \text{adopted}, \text{ballot_num}, pvals \rangle$
 - in active mode spawns commanders for each of the proposal it holds

How a leader enforces C2

- Suppose λ learns that a majority of acceptors has adopted its ballot b ($\langle \text{adopted}, b, pvals \rangle$)
 - CASE 1:** if for some slot s there is no value in $pvals$, then it is impossible that a permanent mapping for a smaller ballot already exists or will ever exist for s : any proposal by λ will satisfy C2

How a leader enforces C2

- Suppose λ learns that a majority of acceptors has adopted its ballot b ($\langle \text{adopted}, b, pvals \rangle$)
 - CASE 2:** let $\langle b', s, p \rangle$ be the pvalue with the maximum ballot number b' for s .
 - by induction, no pvalue other than p could have been chosen for s when $\langle b', s, p \rangle$ was proposed
 - since a majority of acceptors has adopted b , no pvalues between b' and b can be chosen
 - by proposing p with ballot b , λ enforces C2

```

process Leader(acceptors, replicas)
  var ballot_num := (0, self()), active = false, proposals := {}
  spawn(Scout(self(), acceptors, ballot_num));
  for ever
    switch receive();
      case <propose, s, p> :
        if  $\nexists p' : \langle s, p' \rangle \in \text{proposals}$  then
          proposals := proposals  $\cup \{ \langle s, p \rangle \}$ 
          if active then
            spawn(Commander(self(), acceptors, replicas,  $\langle \text{ballot\_num}, s, p \rangle$ ));
          end if
        end if
      end case
      case <adopted, ballot_num, pvals> :
        proposals = proposals  $\oplus \text{pmax}(pvals)$ 
         $\forall \langle s, p \rangle \in \text{proposals} : \text{spawn}(\text{Commander}(\text{self}(), \text{acceptors}, \text{replicas}, \langle \text{ballot\_num}, s, p \rangle));$ 
        active := true
      end case
      case <preempted, r',  $\lambda'$ > :
        if  $(r', \lambda') > \text{ballot\_num}$  then
          active := false;
          ballot_num :=  $(r' + 1, \text{self}());$ 
          spawn(Scout(self(), acceptors, ballot_num));
        end if
    end switch
  end for
end process
    
```



Leader

$$\begin{aligned}
 x \oplus y &\equiv \{ \langle s, p \rangle \mid \langle s, p \rangle \in y \vee \\
 &\quad (\langle s, p \rangle \in x \wedge \nexists p' : \langle s, p' \rangle \in y) \} \\
 \text{pmax}(pvals) &\equiv \{ \langle s, p \rangle \mid \exists b : \langle b, s, p \rangle \in pvals \wedge \\
 &\quad \forall b', p' : \langle b', s, p' \rangle \in pvals \Rightarrow b' \leq b \}
 \end{aligned}$$

```

        end case
      end switch
    end for
  end process
    
```